

SleepScale: Runtime Joint Speed Scaling and Sleep States Management for Power Efficient Data Centers

Yanpei Liu[†] Stark C. Draper^{*} Nam Sung Kim[†]

[†]University of Wisconsin Madison ^{*}University of Toronto

yliu73@wisc.edu stark.draper@utoronto.ca nskim3@wisc.edu

Abstract

Power consumption in data centers has been growing significantly in recent years. To reduce power, servers are being equipped with increasingly sophisticated power management mechanisms. Different mechanisms offer dramatically different trade-offs between power savings and performance penalties. Considering the complexity, variety, and temporally-varying nature of the applications hosted in a typical data center, intelligently determining which power management policy to use and when is a complicated task.

In this paper we analyze a system model featuring both performance scaling and low-power states. We reveal the interplay between performance scaling and low-power states via intensive simulation and analytic verification. Based on the observations, we present SleepScale, a runtime power management tool designed to efficiently exploit existing power control mechanisms. At run time, SleepScale characterizes power consumption and quality-of-service (QoS) for each low-power state and frequency setting, and selects the best policy for a given QoS constraint. We evaluate SleepScale using workload traces from data centers and achieve significant power savings relative to conventional power management strategies.

1. Introduction

As the workloads allocated to data centers increase so does the economic and environmental footprint of these processing clusters. U.S. data center electricity consumption grew by roughly 36% between 2005 and 2010 to about 2% of domestic usage in 2010 (this is about 77,728,000 MW · hr/year, or 12 billion U.S. dollars) [21]. Further, power usage by U.S. data centers is doubling every five years [1]. These numbers give material economic, societal, and environmental reasons for improving power efficiency.

Much power consumed by data centers does not go toward computation. While configured to peak service demand, data centers regularly operate at lighter load levels, and are often only 20 – 30% utilized [3]. Even in idle mode a server continues to draw roughly 60% of the power of a busy server [3, 17]. Although some data centers (e.g. in Google) recently report low power usage effective (PUE), benefiting from the cooling using cold water directly from rivers or seas, around 40% of costs are still associated with power consumption in other commercial and governmental data centers [11]. Thus, it is

important for data centers to maximize their power efficiency and at the same time to guarantee a certain level of QoS.

To maximize power efficiency, commercial processors and platforms support various power management mechanisms such as dynamic voltage/frequency scaling (DVFS) and low-power (sleep) states. However, it is the power management policies (i.e., the choice of operating voltage/frequency and the determination of when to enter which low-power state) that ultimately govern power efficiency for given workload under target QoS constraints. Thus in this paper, the word “policy” refers to some combination of power control methods such as processing speed and low-power state settings.

To develop effective policies we consider a system model that takes into account important QoS constraints and processor, platform, and workload characteristics (Section 3). Then, to develop engineering insight, we study the effectiveness of various policies for the aforementioned system model under idealized setting (Section 4). In contrast to some previous studies, we demonstrate that there is not a single optimal policy; the optimum policy heavily depends on QoS constraints, as well as on the characteristics of processor, platform, and workload. Through intensive simulation and analytic study, we observe that (1) there exists an optimal joint choice of frequency setting and low-power state; (2) the best low-power state depends on the performance constraint at low utilization; (3) the best low-power state also depends on the job size; (4) the delay to enter a low-power state should be jointly determined with frequency and (5) service time dependency on CPU frequency matters.

Based on our analyses and engineering insights, we develop a power management tool, *SleepScale*, that selects the best combination of frequency and low-power states at runtime. It predicts the characteristics of given workload and determines the optimal power management policy with low overhead. While we note that accurately predicting the workload characteristics of a computing system (i.e., the utilization, service time and inter-arrival time distributions) is the key to effective runtime power management, we also demonstrate that simple prediction techniques can offer sufficient accuracy for real-world utilization traces.

The paper is outlined as follows. We summarize related work in Section 2. Power and operational models are introduced in Section 3. Section 4 presents engineering lessons for an idealized workload. SleepScale is introduced in Section 5

State	Operation & Characteristics
$C0_{(a)}$	Operating active state: there is work to do, voltage & frequency setting adjusted dynamically by DVFS
$C0_{(i)}$	Operating idle state: there is no work to do, voltage & frequency held constant at last DVFS setting
C1	Halt state: clock stops
C3	Sleep state: cache flushed, architectural state maintained, clock stopped
C6	Deep sleep state: architectural state saved to RAM, voltage set to zero

Table 1: CPU power states [17].

and evaluated in Section 6. We conclude in Section 7.

2. Background and Related Work

Modern computer systems are equipped with two classes of power management mechanisms to reduce power consumption: performance scaling and component deactivation. DVFS is a widely used performance scaling technique where one reduces processing speed and voltage in step with reduced utilization. In component deactivation one puts the CPU and other peripherals into a low-power sleep state. Table 1 shows typical power states for the Intel Xeon/Atom family of processors. We now summarize some related work on both types of mechanisms.

Power management via performance scaling. Performance scaling such as DVFS has been widely used for computer systems to provide substantial power savings by varying voltage and frequency [15, 20, 30]. In [9] the problem of optimal speed scaling is formulated as a stochastic dynamic problem and a numerical solution is derived. In [2] the speed scaling problem is studied via an algorithmic approach and also jointly considered with scheduling. Recently, in [5], performance scaling in memory systems is considered. In [4] methods are developed to scale the speed of memory systems and processors in a joint manner.

Power management via sleep states. Various low-power states are designed for halting processors when they are in idle. In [25] the authors propose a method for eliminating idle power in servers by quickly transitioning between a high-performance active state and a single low-power state. The development is based on queuing theory [14]. Recent advances in the $M/G/k$ queue with setup costs [7] extend the case to the multi-server scenario: the impact of data center size is studied in [6] and power allocation in server farms is studied in [8]. In a slightly different vein [27] takes a stochastic optimization approach, optimizing time average system performance using optimization theory. However due to distinct wake-up penalties, determining when to enter what low-power state is a complicated task. In [23] the authors warn of potential problems of using these low-power states and suggest “guarded” mechanisms to avoid negative power savings.

Joint approaches and SleepScale. Recently, it is suggested that halting when the system is idle and using a static rate when busy perform almost as good as an optimal speed scaling mechanism [32]. However the low-power state model in [32] is rather limited and only one single low-power state is considered. In [25] the DVFS mechanism is considered in separation from sleep states. Other approaches have also been proposed, for instance a real-time server probing mechanism is proposed in [34]. A more static approach based on workload profiling is proposed in [24].

Nevertheless, much is still not known about which sleep state/speeds are useful when, and how decisions should be made with low overhead to toggle between power saving policies. SleepScale takes a unified approach, modeling a general system with speed scaling and many low-power states having different characteristics. It seeks to develop a family of policies that jointly manage the setting of the operating frequency and the choice of which sleep state to enter, and more importantly, how these decisions should be made online with low overhead.

3. System Model

We now present a system model that accounts for both DVFS and low-power states. We later use the model to study performance and power consumption under various workloads.

3.1. Power model

We discuss how we model (i) the power consumed by the processor, (ii) the power consumed by peripheral (non-CPU) components such as DRAM, hard disk drive (HDD), network interface card (NIC), and (iii) the latencies involved in transitioning between power states.

First, consider processor power. A CPU in the active $C0_{(a)}$ state (and also in the idling $C0_{(i)}$ state) consumes dynamic power. Power consumption will be proportional to $V^2 f$, where V is the supply voltage and $f \in [0, 1]$ is the DVFS clock frequency scaling factor. In our study, we choose a linear DVFS scenario, where both voltage and frequency are scaled linearly. This assumption falls within the scope of some existing processors (see the datasheet in [18]). Dynamic power consumption in states $C0_{(a)}$ and $C0_{(i)}$ will therefore scale cubically with frequency. In the sequel we consider only the frequency parameter f and assume V to be proportional to f . In sleep state C1 the clock signal is gated. Thus, only leakage power is consumed. Platform components also have low-power states and each supports a subset of the CPU states. Table 3 lists the platform power states and the CPU state (or states) that each supports.

The power consumption of the entire system is the sum of CPU power and platform power. In the following we use the term “state” to encompass both CPU and platform state and denote the combined state by concatenating their notations, e.g., $C0_{(i)}S0_{(i)}$. Table 2 tabulates power consumption num-

Components	Operating	Idle	Sleep	Deep sleep	Deeper Sleep
CPU×1 [17]	130V ² <i>f</i> W (<i>C0</i> _(a))	75V ² <i>f</i> W (<i>C0</i> _(i))	47V ² W (<i>C1</i>)	22 W (<i>C3</i>)	15 W (<i>C6</i>)
Chipset×1 [17]	7.8 W	7.8 W	7.8 W	7.8 W	7.8 W
RAM×6 [17]	23.1 W	10.4 W	10.4 W	10.4 W	3.0 W
HDD×1 [29]	6.2 W	4.6 W	4.6 W	4.6 W	0.8 W
NIC×1 [19]	2.9 W	1.7 W	1.7 W	1.7 W	0.5 W
Fan×1 [25]	10 W	1 W	1 W	1 W	0 W
PSU×1 [25]	70 W	35 W	35 W	35 W	1 W
Platform total	120 W (<i>S0</i> _(a))	60.5 W (<i>S0</i> _(i))	60.5 W (<i>S0</i> _(i))	60.5 W (<i>S0</i> _(i))	13.1 W (<i>S3</i>)

Table 2: Power consumption for different components of a system.

State	Operation & supported CPU state(s)
<i>S0</i> _(a)	Active state: associated with <i>C0</i> _(a) only
<i>S0</i> _(i)	Idle state: associated with other CPU states
<i>S3</i>	Sleep: RAM powered, associated with <i>C6</i> only

Table 3: Platform power states [16].

	<i>C0</i> _(a)	<i>C0</i> _(i)	<i>C1</i>	<i>C3</i>	<i>C6</i>
<i>S0</i> _(a)	0 s	–	–	–	–
<i>S0</i> _(i)	–	0 s	1 – 10 μs	10 – 100 μs	0.1 – 1 ms
<i>S3</i>	–	–	–	–	1 – 10 s

Table 4: Average wake-up latencies [12, 25].

bers for the Xeon family of CPUs and associated platform components. As an example, the power consumption in state *C0*_(i)*S0*_(i) is $75V^2f + 52.7$ W.

Table 4 summarizes the delay incurred by the various possible states in returning to active operation. We note that while it is possible to consider a platform shut-down in which the entire system is turned off, the wake-up latency incurred will be enormous, and thus should be considered at a coarser time granularity than is the focus of this work.

3.2. Operation model

We assume jobs arrive at the system according to some random process with rate λ and are served based on the first-come-first-serve (FCFS) order. The server is equipped with a DVFS mechanism, which affect the service time. In active operation the clock frequency can be scaled by a factor $f \in [0, 1]$ and the time it takes to process each job is scaled correspondingly. For CPU-bound jobs, the resulting (scaled) service times are modeled as a random variable with mean $\frac{1}{\mu f}$ where μ is the max service rate. Setting the frequency to the maximum $f = 1$ yields maximum processing speed $\frac{1}{\mu}$ and setting $f = 0$ stops the server from processing jobs, i.e., the server is in a clock-gated mode. For memory-bound jobs, the service time is modeled as independent of frequency, thus with mean $\frac{1}{\mu}$. Finally, we note that the “utilization” factor $\rho = \frac{\lambda}{\mu}$ is the expected fraction of time the server has jobs to process.

As discussed earlier, in active state *C0*_(a)*S0*_(a) power varies

cubically in f ; hence the power is P_0f^3 for some maximum power P_0 . The system also has n low-power states indexed by i , $1 \leq i \leq n$. Each time the server becomes idle it enters a sequence of low-power states, staying in each for a pre-set amount of time. The server enters state i some τ_i seconds after its queue empties. Naturally, $\tau_1 < \tau_2 < \dots < \tau_n$. Formally, the i th low-power state is characterized by the three-tuple (P_i, τ_i, w_i) where:

- P_i is the power consumed in state i ,
- τ_i is the time at which the server enters state i , measured from the time the queue empties, and
- w_i is the average wake-up latency that the system incurs to return to the active state *C0*_(a)*S0*_(a) from state i .

A job arrival interrupts the low-power state and wakes up the system from its current low-power state. We assume a job arrival immediately triggers the wake-up process, during which no job can be served. Deeper sleep states consume less power but take longer to wake up from so $P_1 > P_2 > \dots > P_n$ but $w_1 < w_2 < \dots < w_n$. For simplicity and conservative evaluation we assume the power consumption during the wake-up is the same as the power consumed in active operation.

The low-power states that we study include *C0*_(i)*S0*_(i), *C1**S0*_(i), *C3**S0*_(i), *C6**S0*_(i) and *C6**S3*. We also analyze sequences of pairs of low-power states such as entering *C0*_(i)*S0*_(i) first then *C6**S3*, and also sequences many more low-power states, e.g., entering *C0*_(i)*S0*_(i), *C1**S0*_(i), *C3**S0*_(i), *C6**S0*_(i) and *C6**S3* in sequence.

4. Workload-dependent Optimal Policy

In this section we study the model of Section 3 under the idealized assumptions of Poisson arrivals and exponentially distributed service times. We present our methodology for evaluating various policy choices and discuss engineering lessons.

4.1. Methodology

To evaluate a candidate policy we generate $N = 10,000$ jobs and evaluate the policy at each possible frequency setting based on our model presented in Section 3. The simulated maximum frequency is $f = 1$ and the minimum is the one that the system is barely stable, i.e., $f = \rho + 0.01$ with step

Algorithm 1 Simulation under $\rho = \frac{\lambda}{\mu}$ and frequency f

- 1: Generate jobs with job size and inter-arrival time sampled from probability distributions with mean $\frac{1}{\mu f}$ and $\frac{1}{\lambda}$ respectively (assuming CPU-bound).
 - 2: **for** job j in 1 to N **do**
 - 3: **if** job j arrives before $j - 1$ departs **then**
 - 4: Active += service time of j .
 - 5: Delay of j = departure time of j - arrival time of j .
 - 6: **else**
 - 7: Active += service time of j + wake-up latency.
 - 8: Idle += arrival time of j - departure time of $j - 1$.
 - 9: Delay of j = service time of j + wake-up latency.
 - 10: **end if**
 - 11: **end for**
 - 12: Compute delay by taking average of all j jobs.
 - 13: Compute power by the ratio of active and idle periods.
-

size of 0.01. (We take such a fine step size only to generate smooth plots, in a real system there would be about 10 distinct frequencies.) For policies that consist of a sequence of low-power states, when the processor queue empties the processor enters the low-power states one after another. The arrival of the next job wakes up the processor, incurring some latency. From our simulations we gather results on mean response time $\mathbb{E}[R]$ and average power $\mathbb{E}[P]$. An example of our simulator is provided in Algorithm 1 for single low-power state with $\tau_1 = 0$. Simulating one policy, i.e., one frequency and low-power state combination takes on average 6.3 ms on an Intel i5 2.6 GHz machine using Matlab. This can take even less time when an optimized code is dedicated to run Algorithm 1.

4.2. Engineering lessons

We first consider $\tau_1 = 0$, i.e., whenever the server completes all jobs in its queue the server immediately enters a low-power state from the active state $C0_{(a)}S0_{(a)}$. Second, to investigate how job size effect the choice of policy, we consider two different sizes: “Google-like” ($\frac{1}{\mu} = 4.2$ ms) and “DNS-like” ($\frac{1}{\mu} = 194$ ms). In Section 5.1, when we use traces from data centers for our analysis, the former is roughly the size of a Google web search job and the latter of a DNS look-up job (cf. Table 5). Unless otherwise specified, jobs are assumed to be CPU-bound (and thus service time scales linearly in frequency).

The average wake-up latencies for the system to return to the active $C0_{(a)}S0_{(a)}$ state from various low-power states are set as follows. To wake-up from $C1S0_{(i)}$ we set the average latency to 10 μ s, from $C3S0_{(i)}$ we set it to 100 μ s, from $C6S0_{(i)}$ we set it to 1 ms, and from $C6S3$ we set it to 1 s. These values fall in the ranges specified in Table 4. We observe that other choices from the range specified do not greatly change the following engineering lessons.

We begin our engineering lesson by introducing Figure 1. In each sub-figure we plot average power consumption $\mathbb{E}[P]$

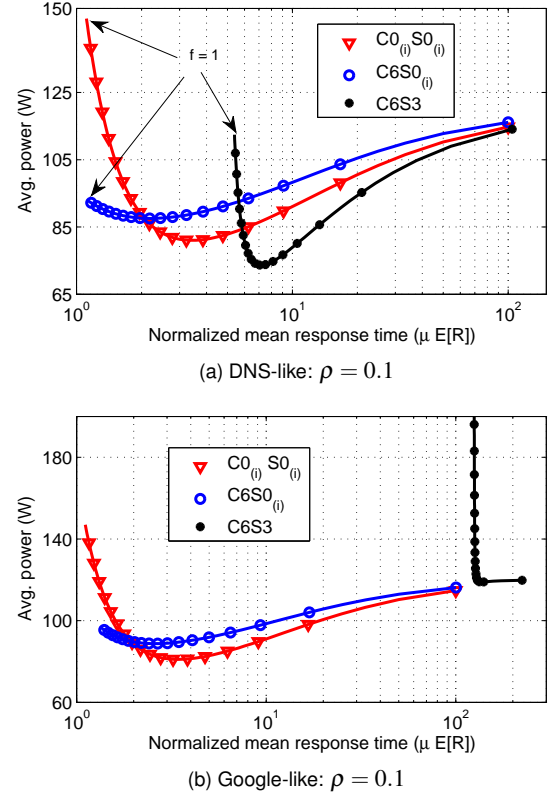


Figure 1: Mean response time $\mathbb{E}[R]$ and average power $\mathbb{E}[P]$ trade-off. Only representative low-power states are shown for the clarity of illustration.

as a function of two parameters. The first parameter is the mean response time $\mathbb{E}[R]$ normalized by the average service time $\frac{1}{\mu}$. We normalize the mean response time in order to compare different workloads (with different average job sizes) directly. The second parameter is the DVFS frequency setting f . Selected choices of f are indicated by the hash marks on each curve. As f is changed the effective service rate μf varies proportionally. The left end of each curve is $f = 1$, corresponding to the fastest processing, and hence smallest response time but most power. The right end of the curve is set by the smallest frequency under which the system is stable, roughly $f = \rho$, where ρ is the utilization. The hash marks are spaced uniformly in increments of 0.05.

We also evaluate systems based on Atom processors with the power numbers from [12]. While simulation results will not be shown, we will discuss our observations.

1) **There exists an optimal joint choice of frequency setting and low-power state.** Figure 1 plots the DNS-like and Google-like workloads results for $\rho = 0.1$. We observe that there is an optimal joint choice of frequency setting and low-power state that yields the minimum power consumption. Without a constraint on the mean response time, the optimal choice corresponds to the bottom of the lowest bowl. For example, in Figure 1-(a) the globally optimum policy uses $C6S3$ and $f = 0.42$ and runs at an average power of 70 W. Setting f

too high leads to worsening efficiency since power increases cubically in frequency. But, setting it too small means that each job takes longer to complete, thereby reducing the possibility of entering a low-power state. The optimal choice – the policy at the bottom of the bowl – strikes a balance between these two effects. Policies such as “race-to-halt” [25] wherein the server runs at maximum frequency $f = 1$ until the queue empties and then immediately enters a single low-power state can consume 50% more power. Such a policy corresponds to the leftmost tip of each curve.

Due to small processor power and relatively large platform power, for Atom processors running DNS-like jobs at low utilizations, it is better to run fast and enter low-power state immediately after the job queue empties.

2) At low utilization, the best low-power state depends on the mean response time budget. In Figures 1 the left-hand side of the plots corresponds to a tight requirement on normalized mean response time. As the constraint is loosened different choices of policies and operating frequency settings become optimal. For example, in Figure 1-(a) under the tightest constraint (when the normalized mean response time $\mu \mathbb{E}[R]$ is required to be in the range $[1, 2]$) policies using $C6S0_{(i)}$ are optimal; under a mid-range constraint policies using $C0_{(i)}S0_{(i)}$ are the best; and under the loosest constraint policies using $C6S3$ prove to be the best. The intuition is as follows.

Consider the range of normalized mean response time for which $C6S0_{(i)}$ is the best in Figure 1-(a). This is the best choice for the tightest response time requirements because setting the frequency high results in faster job completion, thereby increasing the opportunities to enter a more aggressive power saving mode such as $C6$. On the other hand, under a looser constraint on response time the frequency can be lowered yielding the cubic decrease in power consumption. In this situation the expected duration of an empty queue is small and thus it becomes too costly to enter states like $C6$ that have high wake-up latencies. Thus $C0_{(i)}S0_{(i)}$ outperforms $C6S0_{(i)}$ in this situation. This observation is also valid for systems of Atom processors running Google-like workload.

3) The best power state depends on the job size. Figure 2 shows the optimal low-power states among other possible ones (not all shown) for servers under high utilization. When heavily utilized the server rarely exits the active operating state $C0_{(a)}S0_{(a)}$. This reduces opportunities to realize power savings by transitioning to a sleep state, thus most power savings must come from DVFS.

However the job size also plays a crucial rule in the optimal low-power state. For the DNS-like workload, policies using $C6S0_{(i)}$ dominate since the wake-up latency of $C6S0_{(i)}$ is negligible compared to average job size (194 ms). But, for Google-like workload the relatively small job size is sensitive to large wake-up latencies, thus policies using $C3S0_{(i)}$ becomes optimal. For the same reason, very aggressive sleep states such as $C6S3$ should not be used for small size jobs or should be used only during extremely long idle period,

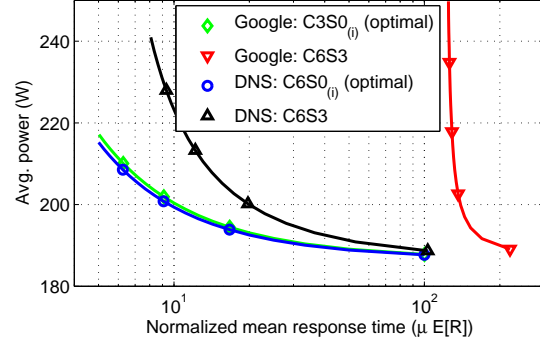


Figure 2: Optimal low-power states for Google and DNS-like workload under high utilization. Other sub-optimal low-power states are not shown for the clarity purpose.

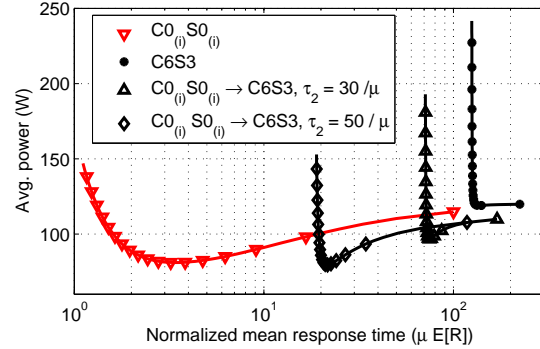


Figure 3: Entering the second low-power state after delays for Google-like workload.

“guarded” by workload prediction techniques [23]. Similar observations can also be made in Atom-based systems.

4) The delay τ_i to enter a low-power state should be jointly determined with frequency. When a server becomes idle, it may wait some amount of time before entering a low-power state to avoid unnecessary wake-up costs [6]. Of course, for some low-power states with very small wake-up latencies, there is no need to delay the entrance as the wake-up incurs negligible penalty. However, for other low-power states with heavy wake-up penalties, it is not immediately clear how long the system should wait.

Figure 3 considers this situation and show what happens if we delay the entrance to state $C6S3$ by various amounts. In “ $C0_{(i)}S0_{(i)} \rightarrow C6S3$ ”, the server first enters $C0_{(i)}S0_{(i)}$ immediately ($\tau_1 = 0$) whenever the job queue empties and will enter $C6S3$ if it idles for some τ_2 seconds (recall the definition of τ_i in Section 3.2). We compare policies using delayed $C6S3$ with ones using immediate $C6S3$ and immediate $C0_{(i)}S0_{(i)}$. We observe that the delay parameter τ_2 interpolates between the immediate $C6S3$ and $C0_{(i)}S0_{(i)}$ curves: setting $\tau_2 = 0$ reduces to immediate $C6S3$ and setting $\tau_2 = \infty$ reduces to immediate $C0_{(i)}S0_{(i)}$. From the plots we observe that by delaying $C6S3$, more power savings can be made at mild mean response time budget (e.g. consider $\mu \mathbb{E}[R] = 20$). Essentially, what we

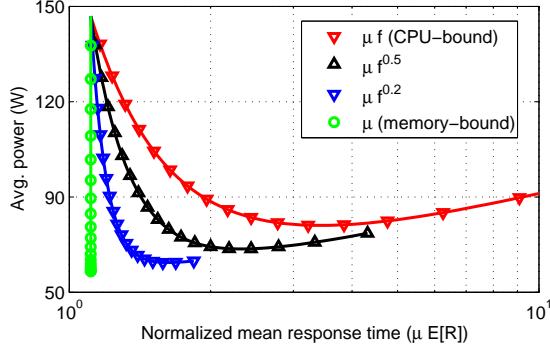


Figure 4: Comparison between different CPU usages for DNS-like workload.

observe is that there is an optimal combination between frequency and entrance delay that minimizes the power under a certain mean response time constraint.

5) Sequential power throttle-back is conservative. It is tempting to concatenate all low-power states i.e., building a system with large number of low-power states and then letting the system enter those states in sequence to derive the most benefit from each. However from our intensive simulation on entering $C0_{(i)}S0_{(i)}$, $C1S0_{(i)}$, $C3S0_{(i)}$, $C6S0_{(i)}$ and $C6S3$ in sequence, we discover that such policies are not often efficient. The reason is that at high utilization the system rarely enters the last state. At low utilization it is a waste of power not to go to the optimal state immediately. Nevertheless, such sequential power reducing policies can be useful when the arrival statistics are unknown.

6) Service time dependency on CPU frequency matters. Now we discuss what happens when the workload is *not* CPU-bound. Recall that for CPU-bound jobs the service rate μ scales linearly with frequency f . For less CPU-bound jobs, the service rate scales sub-linearly and in the extreme case (memory-bound) service rate is insensitive to frequency (as job completion time is dominated by memory access time rather than processing time). In Figure 4 we plot a DNS-like workload at low utilization for service rate varying as different functions of clock frequency. It can be observed that the optimal choice of frequency depends on the scaling. For memory-bound jobs, the optimal speed is the lowest speed.

4.3. Analytic results

In fact, under the Poisson job arrival process and exponential service time distribution we can derive in closed-form the average power $\mathbb{E}[P]$ and mean response time $\mathbb{E}[R]$ for the system model described in Section 3. The results obtained from the closed-form expressions match those presented in Figure 1. Note that constructing Figure 1 using closed-form expression does not involve simulations described in Section 4.1. The closed-form solution can also be derived when the service time follows general distributions, i.e., not limited to exponential distribution [14]. However for both general arrival process

and service time distribution (which SleepScale assumes, as will be discussed next), no closed-form solution exists to the best of our knowledge. We include these theoretical results in the Appendix for reference.

5. SleepScale

In this section we present SleepScale, a runtime power management algorithm. It consists of a policy manager and a runtime predictor. First, in Section 5.1 we detail the policy manager that selects the optimum policy as a function of workload statistics. Second, in Section 5.2 we describe our runtime predictor of the statistical characteristics of the workload such as utilization, and arrival rate and service time distributions. This allows SleepScale to select the best policy in an online manner.

5.1. Policy manager

In this section we describe the policy manager. The manager takes a statistical description of the current workload as input and determines the best policy.

5.1.1. Policy characterization and selection. The policy manager bases its determination of the best policy on the empirically observed distributions of recent arrivals and service times, collected from the server at runtime. The observed recent arrivals and service times can be arbitrary statistics, not limited to Poisson process and exponential distributions. Given the collected statistics it characterizes the power-performance trade-off of each low-power state at a range of frequency settings. It does this by simulating the queuing process as described in Algorithm 1. The optimal policy is the policy that minimizes power consumption while meeting a target QoS constraint. As noted in Section 4.1, simulating a single policy takes only 6 ms. Considering the finite number of frequencies and low-power states, the overhead of simulating all policies is thus negligible compared to the policy updating period (which will be measured in minutes, as discussed later).

Our QoS constraint is determined by a baseline system. Our baseline is motivated by the fact that data centers are often provisioned to meet a QoS target for some peak demand, often specified in an SLA. To meet SLA commitments during periods of peak demand, the data center should be running full out, i.e., at maximum frequency $f = 1$ without using a low-power state. In contrast, at lower loads there is slack in meeting the QoS which can be exploited to reduce operating costs (e.g., power) as much as possible.

We parameterize the target peak demand through a peak design utilization ρ_b . To understand the baseline QoS, consider Figure 5. This figure plots the power-delay trade-off for the Google-like workload running with low-power state $C0_{(i)}S0_{(i)}$ at different frequency settings and under different utilizations $\rho < \rho_b$. In the plot the baseline QoS throughput constraint is indicated by the vertical bar. This vertical bar indicates that the allowable normalized response time is 5 when $\rho_b = 0.8$, calculated (under the idealized model) as

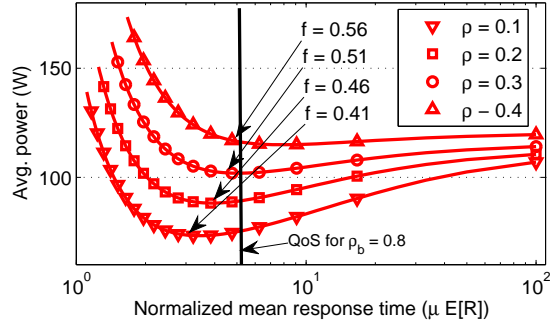


Figure 5: Average power/performance trade-off for Google-like workload.

Workload	Inter-arrival Mean	Inter-arrival C_v	Service Mean	Service C_v
DNS	1.1 s	1.1	194 ms	1.0
Mail	206 ms	1.9	92 ms	3.6
Google	319 μ s	1.2	4.2 ms	1.1

Table 5: Different workload types from [26] (partial list).

$\mu \mathbb{E}[R] = \frac{1}{(1-\rho_b)} = \frac{1}{1-0.8} = 5$. In general as the utilization ρ is increased from $\rho = 0$ to $\rho = \rho_b$ the curves shift up, meaning that a higher frequency setting is required to maintain the QoS throughput constraint. For instance, for $\rho = 0.4$, which is strictly less than $\rho_b = 0.8$, to minimize average power one must set $f = 0.56$ and the system will operate exactly at the requisite QoS. However, for even lower utilizations, e.g., $\rho = 0.1$, one operates at the lowest average power by setting $f = 0.41$. This is the global minimum for this utilization and the normalized mean response time achieved is about 3, which exceeds the QoS requirement. Thus, one can sometimes exceed the baseline QoS while minimizing power consumption.

5.1.2. Results and observations. We now present the results of our policy characterization according to the workload statistics used by the BigHouse [26] simulator. The BigHouse simulator stores statistics of inter-arrival and service times accumulated from long-term observation of live traffic traces for various real-world workloads. Table 5 lists summary statistics (inter-arrival and service time mean and coefficient of variation (C_v)) of three workloads in the BigHouse simulator.

Figure 6 shows the optimum policy as a function of workload, utilization, and QoS constraint (ρ_b). As an example, Figure 6-(a) plots the optimum policy for the DNS-like workload. Operating frequency is indicated on the vertical axis, utilization on the horizontal axis, and the optimum choice of low-power state by the hash marks is indicated in the legend. Depending on the utilization two different policies become optimal. At low utilization $C0_{(i)}S0_{(i)}$ is optimal and at high utilization $C6S0_{(i)}$ is optimal. Two pairs of curves are plotted. The upper two curves are plotted for a baseline QoS constraint set by $\rho_b = 0.6$, the bottom two for $\rho_b = 0.8$. Note that the constraint under $\rho_b = 0.6$ is *tighter* than the one under 0.8.

In each pair of curves one is dashed, meaning it is the policy choice based on the statistics of the BigHouse simulator for that particular workload. The other curve is solid, meaning that it is the optimum policy choice computed by the model considered in Section 4 (i.e., Poisson job arrivals and exponential service times) with the same mean inter-arrival and service time as its paired BigHouse curve.

The plots in the second row of Figure 6 are defined analogously except that the QoS is measured in terms of the 95th percentile response time, rather than in terms of normalized mean response time. Each sub-plot in the second row shares the same workload as the plot directly above. We summarize key observations below.

1) **There is no “one-size-fits-all” policy.** Different workloads and different utilizations require different policies. Almost all low-power states are useful for some set of operating conditions. Thus, relying on a single low-power state when designing power-efficient architectures can be a poor choice.

2) **The idealized model is sometimes good, but often one needs to use more realistic models.** Recall that the model of Section 4 is limited to idealized Poisson job arrivals and exponential service times. These are analytically tractable distributions and when they closely match the actual workload statistics policies based on those results can perform almost as well as the policies based on actual statistics.

We also note that the discrepancy between the policy results based on the idealized model and the BigHouse model is different for two performance constraints; cf. Figures 6-(c) and -(d). This is due to the fact that while the mean response time is only concerned with the mean, the 95th percentile response time constraint is concerned with the tail behavior of the response time distribution, which depends critically on the variation in job size and inter-arrival time.

3) **Often the idealized model computes the best choice of low-power state, but not the frequency setting.** Often for a given utilization the computed optimal low-power state is the same but the frequency setting computed by the idealized model is lower than the one computed by BigHouse. If there is a way to adjust the frequency in runtime, one can rely simply on the idealized model without simulation to compute the optimal policy. We leave this as a part of our future work.

4) **The bump in low utilization region indicates that the policy is exceeding its QoS constraint.** At low utilization, the optimal frequency curve for $C0_{(i)}S0_{(i)}$ often has a concave shape. The consistency of this shape can be explained by referring back to Figure 5. As was observed earlier, at low utilizations the QoS constraint can be exceeded while reaching the global minimum power (optimized across all frequencies). In terms of Figure 5 this means that the systems is operating strictly to the left of the vertical bar. Since the same model underlies the policy optimization based on the idealized and BigHouse models, the global power minimum will be the same for both. For this reason the BigHouse and idealized curves can overlap at low enough utilizations.

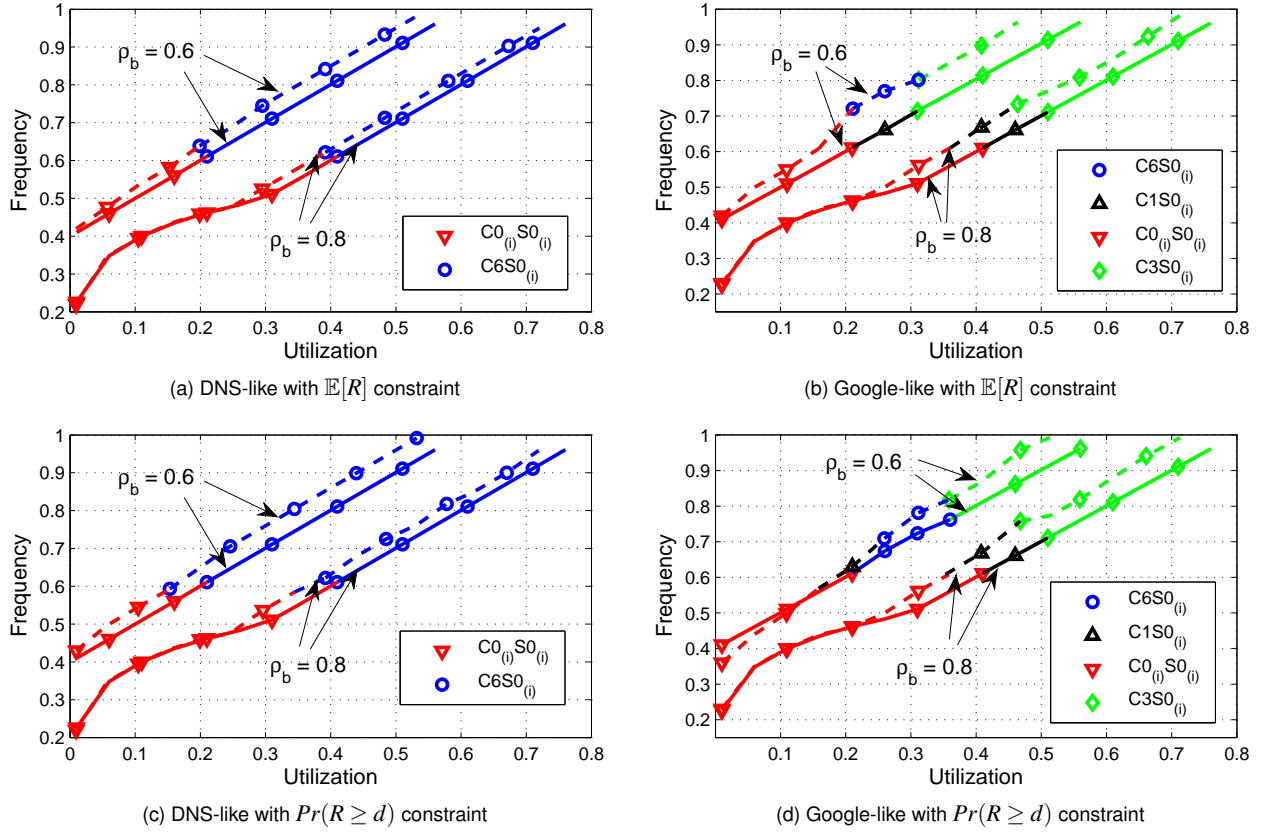


Figure 6: Policy selection for DNS and Google-like workloads. Each curve plots the optimal pairing of frequency setting and low-power states as a function of utilization ρ , and parameterized by QoS constraint ρ_b . The different markers (see legend) indicate which low-power state is optimal at each utilization. The solid lines represent what an idealized model computes and dashed lines represent the results using real-world workload statistics.

As the utilization increases, the optimal frequency setting also increases to keep the power at the global minimum. However beyond some utilization level (roughly $\rho = 0.3$ in Figure 5) the lowest power will no longer continue to meet the target performance constraint. At that level the frequency needs to increase more quickly to continue to meet the constraint. Above that utilization levels the $C0_{(i)}S0_{(i)}$ curves in Figure 6 transition into their respective linear regimes. Since the $\rho_b = 0.6$ constraint is *tighter* than the $\rho_b = 0.8$ constraint and so in curves with arrow “ $\rho_b = 0.6$ ” in Figure 6 we see no evidence of the bump.

5.2. Runtime predictor

The runtime predictor works epoch by epoch, predicting for the current epoch two important aspects of the workload based on its history: 1) inter-arrival time and service time statistics and 2) utilization. Each epoch is a T minutes long ($T \geq 1$) time period. The prediction is fed into the policy manager. The policy is updated at the beginning of each epoch and is held constant throughout the epoch.

5.2.1. Distribution prediction. The inter-arrival time and service time distributions are predicted based upon jobs events

Algorithm 2 Pseudo-code for utilization predictor

```

1: Initialize history depth hist.
2: Initialize a weight vector of size  $p = \text{hist}$ :  $\mathbf{v} = \{v(1), \dots, v(p)\}$ . Set each of the  $p$  entries to  $1/p$ .
3: while prediction for  $\rho(t)$  do
4:   { // Predict the utilization at time  $t$ ,  $\rho'(t)$  using LMS: }
5:   Predict  $\rho'(t) = \min[\sum_{i=1}^p v(i)\rho(t-i), 1]$  from the past  $p$  utilization values.
6:   Compute  $\text{error} = |\rho(t) - \rho'(t)|$ .
7:   Update weight  $\mathbf{v}$  based on  $\text{error}$  and  $\rho(t-1:t-p)$ .
8:   if  $\text{error}$  is larger than some adaptive threshold then
9:     { // CUSUM test: }
10:    Reset  $p = 1$  and set  $v(1) = \text{sum}(\mathbf{v})$ .
11:   else
12:     Grow  $p$ ,  $p = \min(p = p + 1, \text{hist})$  and set  $v(i) = \text{sum}(\mathbf{v})/p$  for all  $1 \leq i \leq p$ .
13:   end if
14:    $t = t + 1$ 
15: end while

```

logged in previous epochs. The logs we collect detail the arrival and service times of each jobs. These logged statistics

are scaled by the predicted utilization (to be discussed later) and fed into the policy manager. Using Algorithm 1 the policy manager then computes the predicted power consumption and QoS of each candidate policy and then selects the policy to use. Note that since we have already obtained the workload log, generating jobs by sampling the distribution (step 1 in Algorithm 1) is not needed. Logging all job events is not necessary either: average behavior from the past several epochs will suffice.

Implicitly the predictor predicts the inter-arrival time and service time distributions based on the past epochs. The motivation for working with the logged arrival and service times is that constructing, maintaining and updating a fine-grained distribution histogram and simulation via sampling is expensive in both time and space. Thus we find it is effective to use logs from previous epochs without explicitly building a distribution. As shown in Section 4.1, it takes only 6.3 *ms* for evaluating one single policy. The policy manager only needs to determine the best policy *once* every epoch. Since, e.g., in the results of Section 6 epochs will be in minutes-length and the determination of the best policy takes less than 1 *s* to compute, the computational overhead is negligible.

5.2.2. Utilization prediction. The distribution predictor predicts the statistics based on the recent epochs. We further enhance such prediction by a fine-grained, minute-by-minute utilization prediction. The workload log gathered in Section 5.2.1 used to simulate the policies is adjusted based on the predicted utilization of the upcoming epoch (the first minute of the epoch, by default): the empirical inter-arrival times between jobs are scaled to match the upcoming predicted utilization.

There is a large body of literature on utilization prediction; much is based on pattern matching [10] across days. As an illustration, in Figure 7 we plot several days worth of minute-granularity utilization traces from academic departmental servers. We observe a periodic daily pattern to the utilization. The abrupt surges observed towards the end of each day in the email store workload, are due to maintenance and back-up services. In contrast to the earlier work based on pattern matching, we study fine-grained minute-by-minute fluctuations in workload behavior. Our approach lets the policy manager react at the processor level to real-time fluctuations in the workload. Our examination will reveal some fundamental aspects of what constitutes good prediction and how the predictor should interact with the selected policy.

In SleepScale we implement three different utilization predictors: a naive-previous predictor, a least-mean-square adaptive filter (LMS) [13], and an LMS filter in conjunction with a cumulative sum change point detection (LMS+CUSUM) [28].

The naive-previous predictor simply uses the utilization in the last minute of the past T -minute epoch as the prediction for the current epoch. This predictor is best suited to track sudden changes in utilization, however it does not effectively predict the stationary behavior of the workload.

The LMS adaptive filter predicts the utilization based on

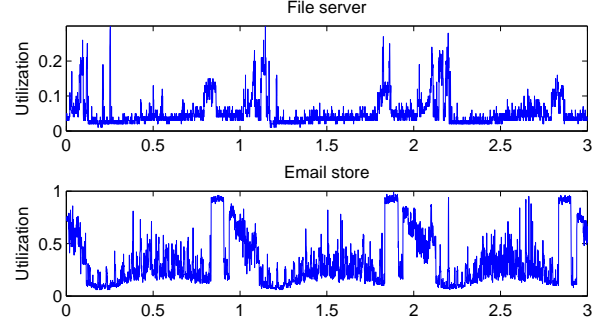


Figure 7: Utilization traces plotted across 3 days for different services, both start at 12 AM of a day. Email store is the host for student, faculty and staff email storage. File server is the host for student files [23, 33].

a weighted combination of the utilizations observed over the past p minutes. The weights are updated every minute based on the prediction error. The LMS adaptive filter outperforms the moving average predictor (which would take the average utilization over the past p minutes) because the weight for each of the past p minutes is chosen adaptively, rather than being fixed to a constant $1/p$. However, like the moving average predictor the LMS filter smooths the data, it does not track abrupt changes well.

As an intermediary between naive-previous predictor and LMS filter, LMS+CUSUM does both tracking and stationary behavior prediction. The pseudo-code of this LMS+CUSUM is given in the Algorithm 2 box. When the CUSUM algorithm detects an abrupt change, the look-back period p in the LMS is reset to 1 (cf. line 10). This resetting drops the smoothing effect of LMS and allows the filter to track the change better. As long as no further abrupt change is detected, p grows until some maximum value is reached (cf. line 12).

5.2.3. Dynamic frequency over-provisioning. The utilization predicted for the first minute of the upcoming T -minute epoch is used to scale the workload log for the entire epoch. The larger T is, the less likely the prediction will be a good one for the entire T -minute epoch. If the predictor overestimates the utilization realized in the epoch, jobs will be processed faster and the queue will tend to empty. However, if the predictor underestimates the realized utilization, the queue will back up, and large delays may result, delays that can propagate into subsequent planning epochs. To control for this in SleepScale we implement the following over-provisioning mechanism.

At the beginning of every T -minute epoch SleepScale computes the average delay incurred by the jobs that were completed in the epoch just past. If the average delay is *below* the delay in the baseline system with utilization ρ_b , i.e., if it is less than $\frac{1}{(1-\rho_b)\mu}$, then the frequency determined by the policy manager is further increased by a factor of α . At first glance this strategy may appear counter-intuitive since one might think it is more natural to over-provision when the past delay has been above (rather than below) the average. However, we

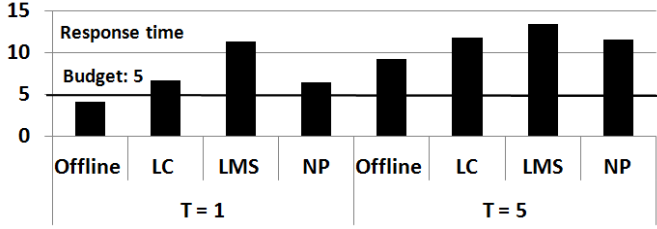


Figure 8: Average response time under different predictors and policy update intervals. No over-provisioning is used (i.e., $\alpha = 0$).

observe that such over-provisioning works best as a “guard band” to buffer against sudden increases in utilization. We illustrate the benefits and the costs of over-provisioning in the next section.

6. Evaluation

To evaluate SleepScale we combine the real-life daily utilization traces discussed in Figure 7 with the BigHouse. We first generate sequences of jobs by sampling the inter-arrival time and service time cumulative distribution functions (CDF) from BigHouse [26]. In systems that serve only a single type of job, the service time distribution is stationary. What varies with utilization is the distribution of inter-arrival times. In our simulated workload traces we then scale the inter-arrival time between generated jobs to match the time-varying utilization of Figure 7. SleepScale uses the job stream as the causal input.

6.1. Under real-world utilization

We first consider a DNS-like server following the email store utilization trace of Figure 7. Across the day the utilization trace covers a large range: from 0.1 to 0.9. We evaluate SleepScale over the period extending from 2 AM to 8 PM as from 8 PM to 2 AM everyday back-up and maintenance operations are scheduled. We select the baseline system to be $\rho_b = 0.8$ resulting in a normalized mean response time budget of $1/(1 - \rho_b) = 5$.

1) **Utilization predictors and policy update interval.** In Figure 8 we study the performance of different predictors: LMS+CUSUM (LC), LMS-only (LMS), naive-previous (NP) and offline (Offline) as well as the policy update periods (T). The offline predictor is a genie-aided predictor where the true utilizations are assumed to be known non-causally in advance. The naive-previous predictor simply uses the utilization in the last minute of the past T -minute epoch as the prediction for the duration of the next T -minute epoch. The LMS+CUSUM and LMS-only predictors are designed with history length $p = 10$. The average response time is measured when running SleepScale with no over-provisioning ($\alpha = 0$).

We note that the more often SleepScale updates its policy (i.e., the smaller the T), the smaller the response time. Since SleepScale selects the best policy based on the predicted utilization, updating the policy fast often helps to mitigate the

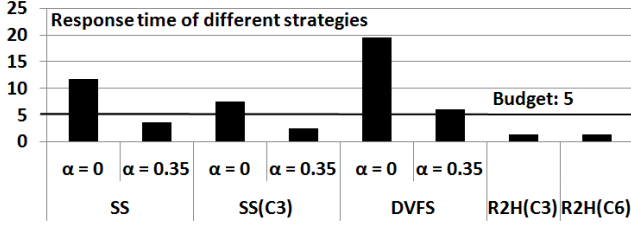
prediction error. We also note that the LMS+CUSUM predictor outperforms the LMS-only predictor since the former can track sudden changes in utilization. However, in some cases this tracking can be detrimental: if a big surge is preceded by a short dip, the tracking predictor may under-estimate the surge since it over-reacts to the dip. For most of the utilization traces that we have, we notice that the accuracy of the naive-previous predictor is often comparable to that of the LMS+CUSUM predictor. The accuracy of these predictors can be further improved by considering the correlation (i.e., repeated daily patterns) across past days.

Finally we note that in Figure 8 the average response time exceeds the allowed budget in all cases when a utilization predictor is used. As mentioned in Section 5.2.3, if the predictor underestimates the realized utilization, the queue will back up, and large delays may result, delays that can propagate into subsequent planning epochs. Thus, we next set the over-provisioning factor $\alpha = 0.35$ and compare the results to the other power management strategies. Recall from Section 5.2.3 that this means the policy manager will increase the frequency by 35% if the delay in the past T -minute epoch is within budget.

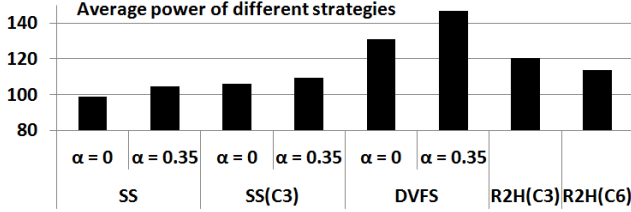
2) **Comparing SleepScale with other strategies.** In Figure 9 we compare SleepScale (SS) with other power control strategies, including a SleepScale method that uses only low-power state $C3S0_{(i)}$ (SS(C3)), a DVFS-only strategy that only uses DVFS and no low-power state (DVFS) and race-to-halt mechanisms using $C3S0_{(i)}$ and $C6S0_{(i)}$ (R2H(C3) and R2H(C6) respectively). Both R2H(C3) and R2H(C6) correspond to the strategy that always operates at the maximum frequency setting ($f = 1$) and transitions into a low power state ($C3S0_{(i)}$ or $C6S0_{(i)}$) immediately upon the queue emptying. All strategies use the LMS+CUSUM predictor with the history length $p = 10$ and are updated every $T = 5$ minutes. We make a number of observations.

SleepScale, when equipped with the frequency over-provisioning, achieves the best power efficiency of all strategies while maintaining the response time within budget. Among others, using DVFS only wastes power as the server is not allowed to enter any low-power state when idling. The race-to-halt mechanism also consumes more power than SleepScale as it sets frequency to the maximum. When SleepScale is set to use only $C3S0_{(i)}$ it also consumes more power as this low-power state is a sub-optimal one. *Our results clearly demonstrate the importance of joint optimization of speed scaling and sleep state selection.*

The properly set over-provisioning reduces response time at the cost of a slight increase in power. This is due to the fact that the extra capacity over-provisioning allocates during periods of low utilization allows the server to accommodate unpredictable surges in utilization. This proves to be essential to meet the mean response time budget. Also running slightly faster does not cost too much power as the server can enter low-power state sooner. (This is not true for DVFS-only strategy



(a) Response time comparison



(b) Average power comparison

Figure 9: Comparing SleepScale with other power control strategies. All strategies are running with LMS+CUSUM predictor with $p = 10$ and are updated every $T = 5$ minutes.

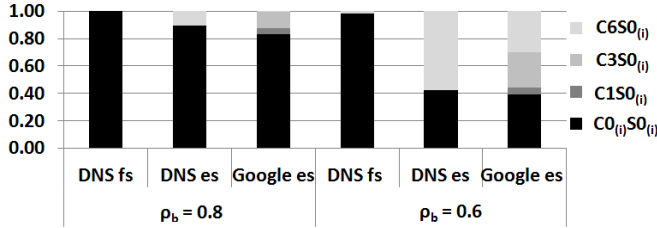


Figure 10: Distribution of optimal low-power states selected by SleepScale. LMS+CUSUM predictor is used with history length $p = 10$, update interval $T = 5$ and over-provisioning $\alpha = 0.35$.

as it has no low-power state to use.)

Finally we comment on the large response time in the DVFS-only strategy. To minimize power under a given response time budget, the optimal solution for the DVFS-only strategy is to set the frequency low enough that the response time just meets the budget. This however consumes all the performance budget thus even a slight utilization miss-prediction can result in large queuing delay. However for SleepScale, as we demonstrate in Figure 5, does not always consume the entire performance budget.

6.2. Distribution of low-power states

In Figure 10 we present the distributions of the low-power states selected by SleepScale for different workloads, base-lines, and utilizations. We report results for the file server (fs) and email store (es) utilizations. We run both DNS and Google-like services with $\rho_b = 0.6$ and $\rho_b = 0.8$.

At low average utilization and when there are few time-varying fluctuations (which is the case for file server), a single low-power state often suffices. For highly time-varying utiliza-

tion traces (such as email store), multiple low-power states are used: $C0(i)S0(i)$ and $C6S0(i)$. Tightening the performance constraint further will lead to deeper sleep states being used more often, as the fast processing required to meet the budget creates more opportunities for entering aggressive power-saving states like $C6S0(i)$. These observations all match our analyses in Section 4.

7. Conclusion and Future Work

In this paper we develop SleepScale, a power management tool to manage data centers for power efficiency while meeting QoS agreements. SleepScale uses queuing-based simulations to determine how best to exploit the low-power states and operating frequency built into modern CPUs. The optimum such policy can be determined at runtime, via online prediction for workload statistics and utilization. We characterize the performance of SleepScale on data center traces. We evaluate SleepScale realizing significant power savings relative to some conventional power management strategies while meeting the same QoS constraints.

More power control knobs are being built into different hardwares in computer systems. However challenges remain in effectively adjusting these knobs in a coordinated manner. One future research direction involves not only the power control mechanisms in CPUs, but also the ones in other system components. We conjecture that high-level queuing models and simulations (rather than fine-grained instruction-level simulation) can help in designing effective power control strategies. Another research direction involves studying SleepScale on multi-core, multi-server systems in order to scale out. The focus is on controlling the overall queuing simulation overhead, although SleepScale can be performed on each core or server independently.

Acknowledgment

The authors would like to thank Aman Chadha (UW-Madison) for early efforts on workload generation, Daniel Wong (USC) for the departmental datacenter utilization traces, Daniel Myres (UW-Madison) for the discussion on queuing theory, Ken Vu, Srinu Ramani (IBM) for their continuous suggestion and support, and the anonymous reviewers for their feedback.

This work is supported in part by NSF grants (CCF-0963834, CCF-0953603, and CNS-1217102) and an AFOSR grant (FA9550-13-1-0138). Nam Sung Kim has a financial interest in AMD.

Appendix

The analytic results extend the work in [22]. The average power consumption for the single-server system with n low-power states described in Section 4 is

$$\mathbb{E}[P] = \frac{1}{\lambda L} \left[\sum_{i=1}^{n-1} P_i (e^{-\lambda \tau_i} - e^{-\lambda \tau_{i+1}}) + P_n e^{-\lambda \tau_n} \right] + P_0 \left(1 - \frac{e^{-\lambda \tau_1}}{\lambda L} \right)$$

where L is defined as

$$L = \frac{\mu f + \mu f \lambda \left[\sum_{i=1}^{n-1} w_i (e^{-\lambda \tau_i} - e^{-\lambda \tau_{i+1}}) + w_n e^{-\lambda \tau_n} \right]}{\lambda (\mu f - \lambda)}.$$

This can be proved using busy period analysis and first principles, see [14]. The mean response time $\mathbb{E}[R]$ is

$$\mathbb{E}[R] = \frac{1}{\mu f - \lambda} + \frac{2\mathbb{E}[D] + \lambda \mathbb{E}[D^2]}{2(1 + \lambda \mathbb{E}[D])},$$

where

$$\mathbb{E}[D^\alpha] = \sum_{i=1}^{n-1} w_i^\alpha (e^{-\lambda \tau_i} - e^{-\lambda \tau_{i+1}}) + w_n^\alpha e^{-\lambda \tau_n}.$$

This can be derived using the result from [31] and some algebraic manipulations. Both $\mathbb{E}[R]$ and $\mathbb{E}[P]$ can be extended to the case where service time is not exponential.

The probability the response time R exceeds deadline d is

$$Pr(R \geq d) = \frac{e^{-(\mu f - \lambda)d} - w_1(\mu f - \lambda)e^{-d/w_1}}{1 - w_1(\mu f - \lambda)}.$$

Note that when $d = 0$ the $Pr(R \geq d) = 1$. When $d = \infty$ the $Pr(R \geq d) = 0$. When $w_1 = 0$ the $Pr(R \geq d) = e^{-(\mu f - \lambda)d}$. When $w_1 = \infty$ the $Pr(R \geq d) = 1$. This result can be proved using a Laplace transform analysis of the response time R .

References

- [1] H. Amur, J. Cipar, V. Gupta, G. R. Ganger, M. A. Kozuch, and K. Schwan, "Robust and flexible power-proportional storage," *Proceedings of the ACM Symposium on Cloud Computing*, pp. 217–228, 2010.
- [2] L. L. H. Andrew, M. Lin, and A. Wierman, "Optimality, fairness and robustness in speed scaling designs," *Proceedings of the ACM SIGMETRICS international conference on Measurement and modeling of computer systems*, pp. 37–48, 2010.
- [3] L. A. Barroso and U. Hözl, "The case for energy-proportional computing," *Computer*, vol. 40, pp. 33–37, Dec. 2007.
- [4] Q. Deng, D. Meisner, A. Bhattacharjee, T. Wenisch, and R. Bianchini, "CoScale: Coordinating CPU and memory system DVFS in server systems," *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 143–154, Dec. 2012.
- [5] Q. Deng, D. Meisner, L. Ramos, T. F. Wenisch, and R. Bianchini, "MemScale: active low-power modes for main memory," *ACM Proceedings of International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, pp. 225–238, 2011.
- [6] A. Gandhi and M. Harchol-Balter, "How data center size impacts the effectiveness of dynamic power management," *Allerton Conference on Communication Control and Computing*, pp. 1164–1169, Sep. 2011.
- [7] A. Gandhi, M. Harchol-Balter, and I. Adan, "Server farms with setup costs," *Performance Evaluation*, vol. 67, pp. 1123–1138, Nov. 2010.
- [8] A. Gandhi, M. Harchol-Balter, R. Das, and C. Lefurgy, "Optimal power allocation in server farms," *Proceedings of the ACM SIGMETRICS Joint International Conference on Measurement and Modeling of Computer Systems*, Jun. 2009.
- [9] J. M. George and J. M. Harrison, "Dynamic control of a queue with adjustable service rate," *Operation Research*, vol. 49, no. 5, pp. 720–731, 2001.
- [10] D. Gmach, J. Rolia, L. Cherkasova, and A. Kemper, "Workload analysis and demand prediction of enterprise data center applications," *Proceedings of the International Symposium on Workload Characterization*, pp. 171–180, 2007.
- [11] A. Greenberg, J. Hamilton, D. A. Maltz, and P. Patel, "The cost of a cloud: research problems in data center networks," *ACM SIGCOMM Computer Communication Review*, vol. 39, pp. 68–73, Dec. 2008.
- [12] M. Guevara, B. Lubin, and B. C. Lee, "Navigating heterogeneous processors with market mechanisms," *to appear in IEEE International Symposium on High Performance Computer Architecture (HPCA)*, Feb. 2013.
- [13] F. Gustafsson, *Adaptive filtering and change detection*. Wiley, Sep. 2000.
- [14] M. Harchol-Balter, *Performance Modeling and Design of Computer Systems: Queueing Theory in Action*. Cambridge University Press, 2013.
- [15] S. Herbert and D. Marculescu, "Analysis of dynamic voltage/frequency scaling in chip-multiprocessors," *ACM/IEEE International Symposium on Low Power Electronics and Design (ISLPED)*, pp. 38–43, Aug. 2007.
- [16] HP, Intel, Microsoft, Phoenix and Toshiba. (2011, Dec.) The ACPI specification. [Online]. Available: <http://www.acpi.info/>
- [17] Intel. (2012, May) Intel Xeon processor E5-1600/E5-2600/E5-4600 product families. [Online]. Available: <http://tinyurl.com/d7ma5nf>
- [18] —. (2013, Jan.) Intel 80200 processor based on Intel XScale microarchitecture datasheet. [Online]. Available: <http://tinyurl.com/15xz77u>
- [19] —. (2013, Mar.) Intel ethernet controller I350 datasheet. [Online]. Available: <http://tinyurl.com/cc8mhvp>
- [20] S. Kaxiras and M. Martonosi, *Computer Architecture Techniques for Power-Efficiency*. Morgan & Claypool, 2008.
- [21] J. Koomey, "Growth in data center electricity use 2005 to 2010," *Analytics Press*, 2011.
- [22] Y. Liu, S. C. Draper, and N. S. Kim, "Queueing theoretic analysis of power-performance tradeoff in power-efficient computing," *IEEE Conference on Information Sciences and Systems (CISS)*, pp. 1–6, Mar. 2013.
- [23] N. Madan, A. Buyuktosunoglu, P. Bose, and M. Annavaram, "A case for guarded power gating for multi-core processors," *IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 291–300, Feb. 2011.
- [24] J. Mars, L. Tang, R. Hundt, K. Shadron, and M. L. Soffa, "Bubble-Up: increasing utilization in modern warehouse scale computers via sensible co-locations," *IEEE/ACM International Symposium on Microarchitecture*, pp. 248–259, 2011.
- [25] D. Meisner, B. T. Gold, and T. F. Wenisch, "PowerNap: eliminating server idle power," *ACM Proceedings of International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, pp. 205–216, Mar. 2009.
- [26] D. Meisner, J. Wu, and T. F. Wenisch, "BigHouse: a simulation infrastructure for data center systems," *Proceedings of the IEEE International Symposium on Performance Analysis of Systems & Software (ISPASS)*, pp. 35–45, 2012.
- [27] M. J. Neely, "Lower power dynamic scheduling for computation systems," Tech. Rep. arXiv:1112.2797, Dec. 2011.
- [28] E. S. Page, "Continuous inspection schemes," *Biometrika Trust*, vol. 41, pp. 100–115, Jun. 1954.
- [29] Seagate. (2012, Dec.) Desktop HDD ST500DM002. [Online]. Available: <http://tinyurl.com/c358blv>
- [30] D. C. Snowdon, S. Ruocco, and G. Heiser, "Power management and dynamic voltage scaling: Myths and facts," *Proceedings of the 2005 Workshop on Power Aware Real-time Computing*, Sep. 2005.
- [31] P. D. Welch, "On the generalized M/G/1 queueing process which the first customer of each busy period receives exceptional service," *Operation Research*, vol. 12, pp. 736–752, Sep. 1964.
- [32] A. Wierman, L. L. H. Andrew, and A. Tang, "Power-aware speed scaling in processor sharing systems," *Performance Evaluation*, vol. 69, pp. 601–622, Dec. 2012.
- [33] D. Wong and M. Annavaram, "KnightShift: scaling the energy proportionality wall through server-level heterogeneity," *Proceedings of IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 119–130, 2012.
- [34] H. Yang, A. Breslow, J. Mars, and L. Tang, "Bubble-Flux: Precise online QoS management for increased utilization in warehouse scale computers," *International Symposium on Computer Architecture*, pp. 607–618, 2013.